

Particle in cell simulations on GPU clusters

Francesco Rossi ¹

¹University of Bologna and INFN, LOASIS affiliate for the summer

AFRD Meeting, May 21st, 2013

Table of Contents

- 1 Introduction to GPUs
- 2 GPU PIC current deposition algorithm
- 3 Multi GPU parallelization
- 4 Benchmarks and simulations



Simulations for laser plasma acceleration

- Particle in cell (PIC) simulations are useful tools for designing and optimizing laser-driven, plasma-based accelerators.
- Such simulations may require a huge amount of computation ($\gtrsim 10^5$ CPU hours).



Motivation: why a GPU PIC code?

2012

Run simulations we need to run today on the most **efficient** parallel architectures available (GPUs) for PICs

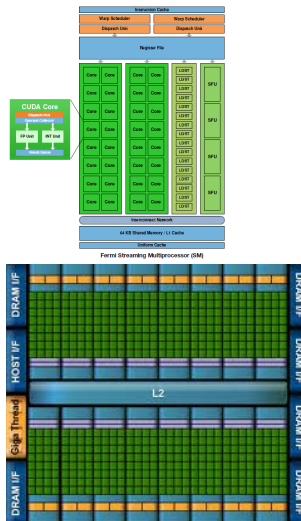
- **Performance** exploiting exposed parallelism
- **Efficiency** from data locality

2016

- Path to exascale computing era → dominated by **manycore** architectures
 - **Prepare PIC algorithms for massively-manycore shared-memory node systems**
 - Bigger subdomains → fundamental for scalability/load balancing
- GPUs roadmap promises **200% performance increase every year** and a half (next generation out this fall)
 - Help to sustain the **computational demand** in LPA community

The NVIDIA CUDA GPU architecture

- On chip: ~15 Multiprocessors, each one:
 - 256 KB register files
 - 16-48 Kb [manual](#) cache: shared memory
 - Issuing instructions
 - executing “warps” of 32 threads in a SIMD fashion
 - divergent branches in a warp cause warp [serialization](#)
- Hides latency keeping many thread warps in flight
- High-bandwidth memory bus (~200 Gb/s) connecting to device RAM
 - Prefers ordered access within a warp
 - Cannot rely on cache: number of [cache bytes-per-thread](#) is several orders of magnitude lower than on CPUs



The latency / throughput dilemma

- Memory **Latency** and **Bandwidth** are often **limiting** performance. Two different strategies:

Single thread optimization

- In scalar processors
- **Reduce latencies**
- Use large caches (per thread)
- Predict branches

Throughput processors: GPUs

- 1 Provide high bandwidth/**throughputs**
- 2 *Saturate it: Tolerate* latencies processing **many threads in parallel**
- 3 Space/energy **saved** removing scalar optimization used for having more computational power



Similar considerations also apply for other instructions, not only memory accesses.

Jasmine

- “*Jasmine*”, a 3D GPU particle in cell code (PIC), featuring:
 - Second order explicit PIC algorithm (FDTD + Boris Pusher) in double precision
 - High order particles shape functions
 - Charge conserving simulations using Esirkepov shape factors
 - 3D multi-GPU simulations with high scalability
 - Dynamic load balancing
 - Moving window
 - Particle trajectory tracking
 - Simulation restart & asynchronous I/O
 - Integrated with a radiation generation computation code

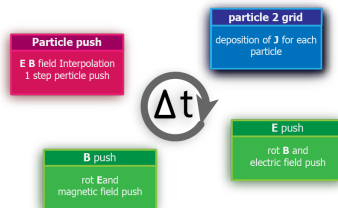


Table of Contents

- 1 Introduction to GPUs
- 2 GPU PIC current deposition algorithm**
- 3 Multi GPU parallelization
- 4 Benchmarks and simulations



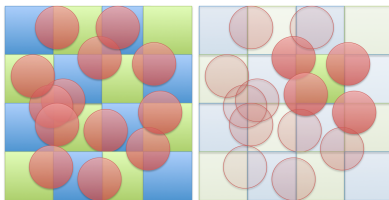
Particle in cell method, algorithmically



- EM PIC: Maxwell+Lorentz+(Vlasov sampling)
- Particle-grid interactions:
 - Force on particles is interpolated (averaging) *from fields grids*
 - Particle current/density is deposited *to grid*
 - scatter operation: a particle adds its density value the cells that it overlaps



PIC Deposition algorithm on massively parallel architectures



- 1 Naive, (1 particle \longleftrightarrow 1 thread) parallelization \rightarrow Race conditions on same memory cells: wrong results
 - Atomic operations or other **synchronization** methods are required
- 2 N particles per cell, particle shape function of total order K (4~27)
 - Density grid data is accessed $K \cdot N_{ppc}$ times
 - It's **worth caching** in GPU multiprocessor's **shared memory**

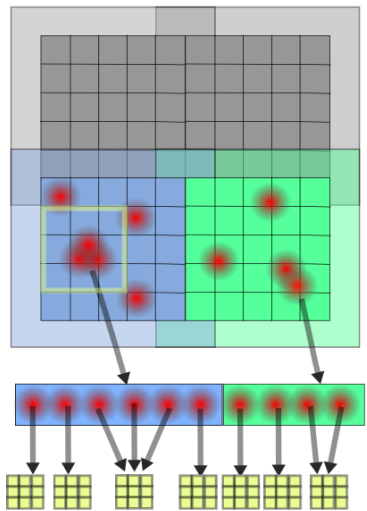


Deposition algorithm without atomics

Algorithm

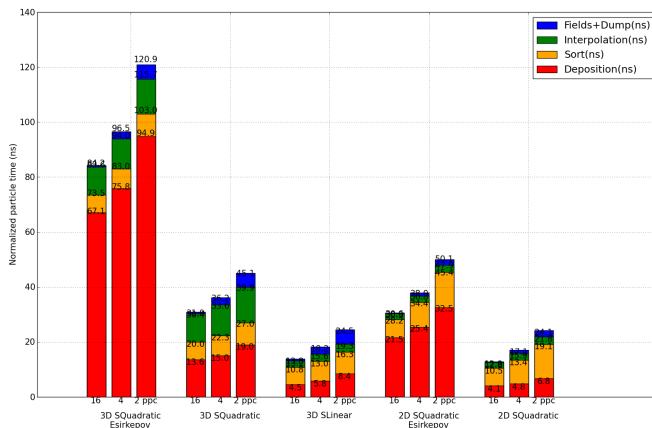
- Sort by cell block and cell
- Assign a CUDA block to a cell block
- Perform a per-block, shared memory, *segmented scan* to compute density sum **for each cell**
- Sum cached copy to global grid

Cell start	1	0	0	0	1	0	0	1
Data	1	2	3	4	0	3	3	1
Pass 1	1	3	5	7	0	3	6	1
Pass 2	1	3	5	7	0	3	6	1
Pass 3	1	3	5	7	0	3	6	1



Performance

Jasmine performance benchmark on NVIDIA Tesla M2075



LASER: $a=7.7$, waist=9.0 μm , fwhm=24.0 fs PLASMA: density= $1.00\text{e}+19$ $1/\text{cm}^3$. Simulation run for $ct = 60\mu\text{m}$, double precision. **Note:** 3D test with Esirkepov method runs stretched grid ▶



Performance gain

- *Jasmine* vs *ALaDyn* (our CPU code), same **exact** simulation.
- Performance of a *single* NVIDIA Fermi GPU equates ~200x BlueGene cores or ~45x IBM SP6 cores.
 - Plus, since subdomains are much larger, load balancing and scalability are much easier
- In the simulation setups shown above (fair resolution), *jasmine* can simulate **~4mm of plasma per day** on a ~24 GPUs cluster



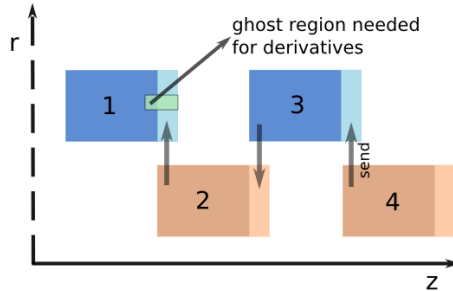
Table of Contents

- 1 Introduction to GPUs
- 2 GPU PIC current deposition algorithm
- 3 Multi GPU parallelization**
- 4 Benchmarks and simulations



Scaling to multiple GPUs: Hiding network transfer

- Exchange:
 - 1 Fields' halos
 - 2 **Particles** leaving subdomain
- Cluster nodes communicate using standard **MPI**



Scaling to multiple GPUs: Hiding network transfer

- Exchange:
 - 1 Fields' halos
 - 2 **Particles** leaving subdomain
- Cluster nodes communicate using standard **MPI**
- Transfer particles **concurrently** with current deposition.
 - Communication can be **hidden** almost completely.

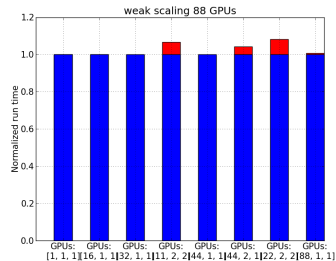
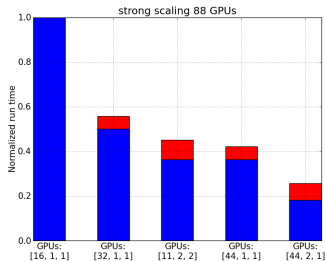
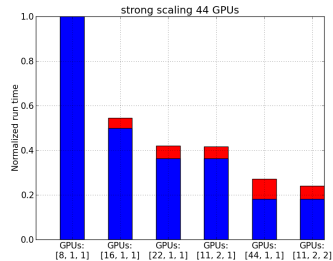
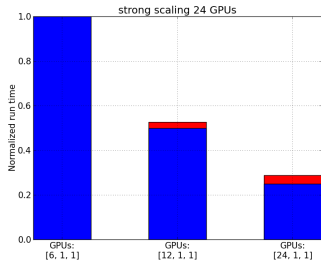


- Scalability test: warm plasma simulation on INFN **APE** cluster @ “ULa Sapienza” and **PLX** machine @ CINECA



Scaling

Warm plasma simulation strong/weak scaling



Simple algorithm for load balancing

- The particle motion easily leads to **inhomogeneous** distribution of the load
- **Shrink** the volume of heavy-loaded nodes:
 - Each few timesteps **select** a subdomain (and its row) and the direction where to shrink
 - Subdomains topology remains intact (vertices conservation)
 - Choice is done trying to **minimize** the cost function:

$$k_1 \times \text{Max}(\text{Load}) / \text{Average}(\text{Load}) + k_2 \times \text{Variance}(\text{Load}) / \text{Average}(\text{Load})$$



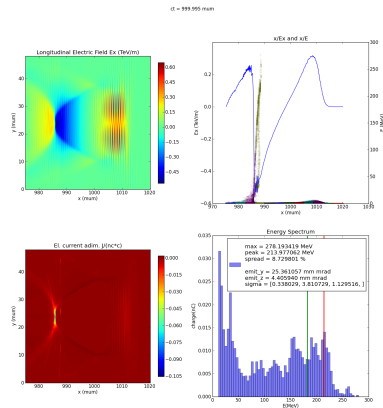
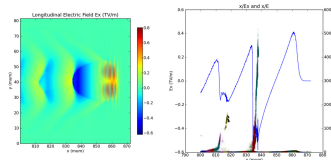
Load balancer test case

Setup (coarse test)

LASER: $a=5.8$, waist=13.2 μm ,
fwhm=24.0 fs

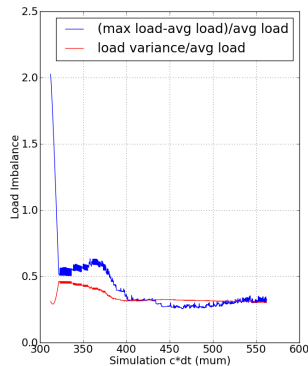
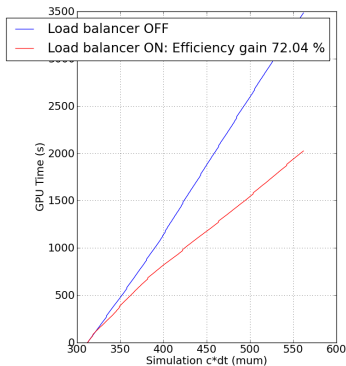
PLASMA: density= $3.80\text{e}+18$ $1/\text{cm}^3$

GRID: $n=[729, 96, 96]$, $dx=['6.25\text{e}-02'$,
' $5.00\text{e}-01'$, ' $5.00\text{e}-01'$] μm



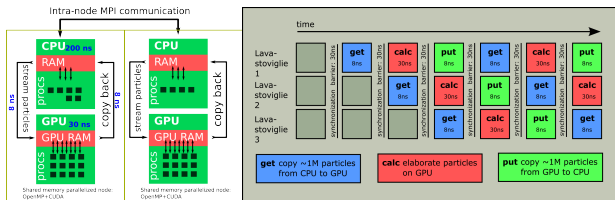
Test case: scaling with load balancing

- With 72 subdomains:



Memory intensive simulations

- Total memory availability represents a constraint for many simulations (for example ion acceleration ones).
 - In a node, GPU memory is often much less than the total host memory available
 - Using host memory to store simulation data make larger simulations **possible** on a cluster of fixed size
- **Asynchronous** stream of particle chunks stored in main CPU memory overlapped with computation using CUDA streams
 - Slower but no longer memory bound to the GPU device memory
 - Currently in testing stage



Implementation note: Meta-programming

- Meta-programming can be used:
 - for writing maintainable code for all particle weighting / numerical schemes
 - **tweak** parameters for optimization of each case
 - implementing **different** numerical schemes using the same **core** algorithms (deposition and interpolation)
- Attempts:
 - 1 **C++** template meta-programming
 - 2 **Python-based** code template engine (code becomes more linear, but non standard)
- Python also used for simulation initial conditions definition, plotting (numpy+matplotlib) and basic automated data-analysis



Table of Contents

- 1 Introduction to GPUs
- 2 GPU PIC current deposition algorithm
- 3 Multi GPU parallelization
- 4 Benchmarks and simulations**



LWFA benchmark simulation

Setup

LASER: $a=2.0$, waist=8.2 μm ,

fwhm=21.0 fs

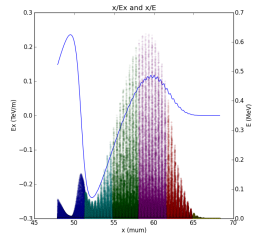
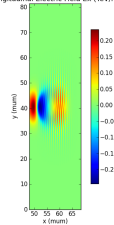
PLASMA: density= $1.38 \times 10^{19} \text{ 1/cm}^3$

GRID: $n=[512, 256, 256]$, $dx=['4.00 \times 10^{-2}'$,

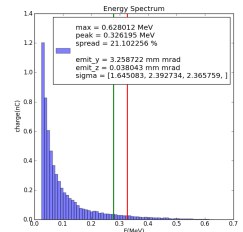
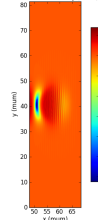
$'3.18 \times 10^{-1}'$, $'3.18 \times 10^{-1}'] \text{ }\mu\text{m}$

ct = 47.951 μm

Longitudinal Electric Field Ex (TeV/m)



El. current adim. $j/(nc^+)$



LWFA benchmark simulation

Setup

Parameters from Paul et al.

*Benchmarking the codes VORPAL,
OSIRIS, and QuickPIC with Laser
Wakefield Acceleration Simulations,*
AIP Conference

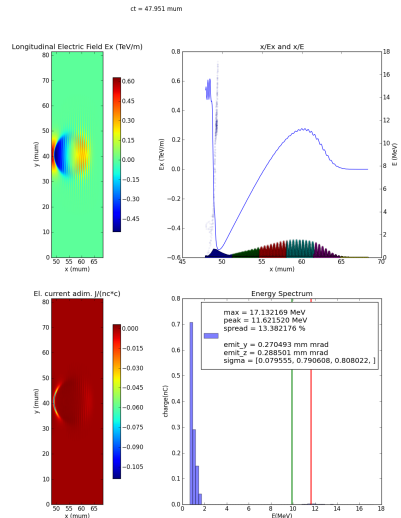
Proceedings;1/22/2008

LASER: $a=4.0$, waist=8.2 μm ,

fwhm=21.0 fs

PLASMA: density= $1.38\text{e}+19$ $1/\text{cm}^3$

GRID: $n=[512, 256, 256]$, $dx=['4.00\text{e}-02'$,
' $3.18\text{e}-01'$, ' $3.18\text{e}-01'$] μm



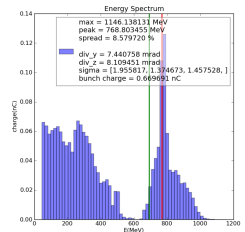
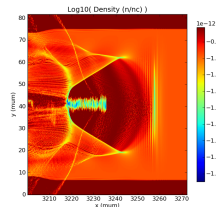
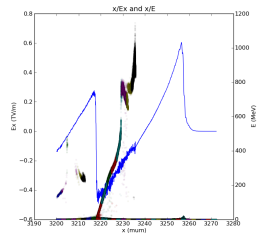
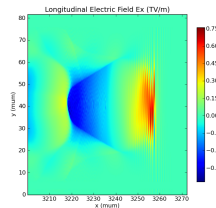
SPARCLab electron acceleration simulation

Setup

LASER: $a=4.9$, waist=15.5 μm ,

fwhm=30.0 fs

PLASMA: density= 3.0×10^{18} $1/\text{cm}^3$



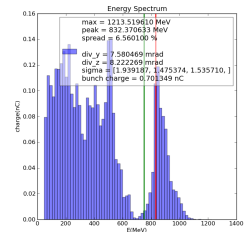
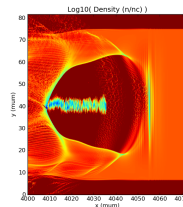
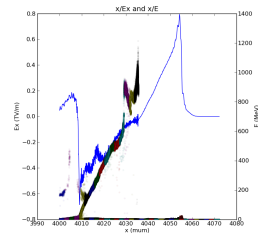
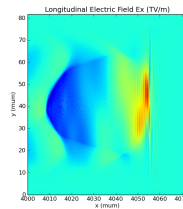
SPARCLab electron acceleration simulation

Setup

LASER: $a=4.9$, waist=15.5 μm ,

fwfm=30.0 fs

PLASMA: density= 3.0×10^{18} $1/\text{cm}^3$



TNSA ion acceleration: Frascati Flame parameters

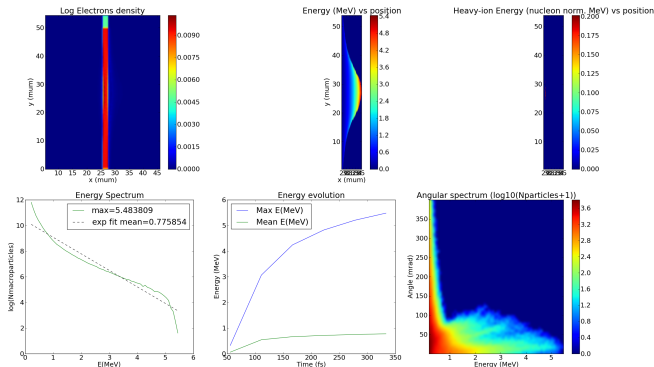
Setup

LASER: $a=7$, waist=10.0 μm , fwhm=30.0 fs

TARGET: 2.0 μm thick

Double layer: aluminium, ($n/n_c = 100$) + back side contaminants layer

$T = 333.332$ fs



TNSA ion acceleration: Nara-like parameters

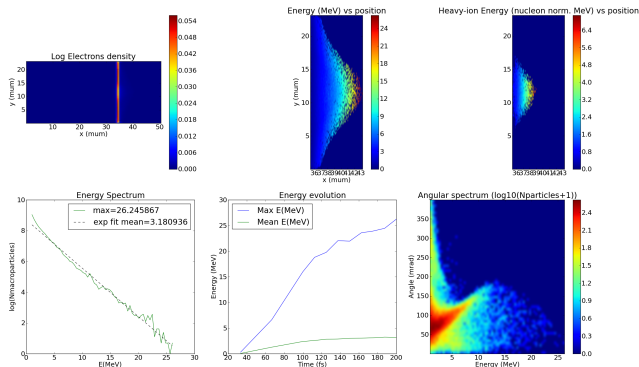
Setup

LASER: $a=22$, waist=3.5 μm , fwhm=40.0 fs

TARGET: 0.8 μm thick

Double layer: aluminium ($n/\text{nc} = 60$) + back side contaminants layer

T = 199.991 fs



Future work

- Performance tuning for Kepler architecture
- Implement more accurate and/or **optimized numerical schemes**.
 - GPUs alone are **not** of enough for satisfying all the computational requirements of the experimental groups (e.g. simulating a **10GeV** electron acceleration stage).
- Lot of work to do!



References

- 1 F.Rossi et al., Towards robust algorithms for current deposition and dynamic load-balancing in a GPU particle in cell code. AIP Conference Proceedings Vol. 1507 Issue 1, p184
- 2 Birdsall,Langdon. Plasma physics via computer simulation.
- 3 Benedetti,Sgattoni,Turchetti,Londrillo. ALaDyn: A High-Accuracy PIC Code for the Maxwell-Vlasov Equations. IEEE Transactions on Plasma Science,36(4),2008.
- 4 Burau,Widera,Hönig,Juckeland,Debus,Kluge,Schramm,Cowan,Sauerbrey,Bussmann. PIConGPU: A Fully Relativistic Particle-in-Cell Code for a GPU Cluster, IEEE Transactions on Plasma Science 38(10)
- 5 Abreu,Fonseca Pereira,Silva. PIC Codes in New Processors:A Full Relativistic PIC Code in CUDA-Enabled Hardware With Direct Visualization.IEEE Transactions on Plasma Science,vol.39,issue 2
- 6 Kong,Huang,Ren,Decyk.Particle-in-cell simulations with charge-conserving current deposition on graphic processing units. J.Comput.Phys.230,4(February2011)
- 7 Sengupta,Harris,Zhang,Owens.Scan primitives for GPU computing.In Graphics Hardware 2007, Aug.2007.
- 8 Hoberock,Bell.Thrust:A Parallel Template Library

